



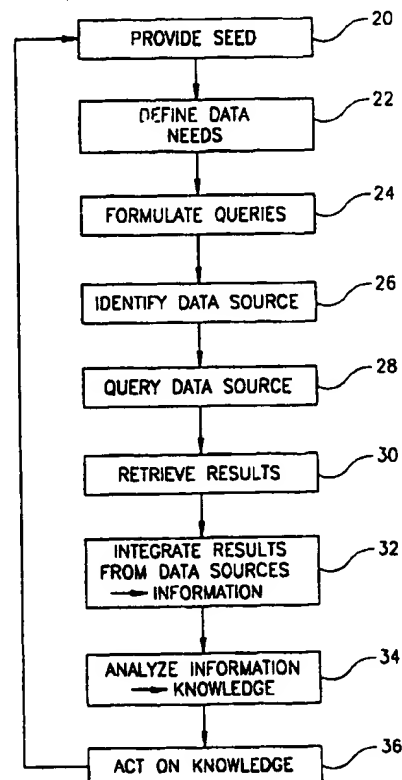
INTERNATIONAL APPLICATION PUBLISHED UNDER THE PATENT COOPERATION TREATY (PCT)

(51) International Patent Classification ⁷ : C12Q 1/68		A2	(11) International Publication Number: WO 00/15847
			(43) International Publication Date: 23 March 2000 (23.03.00)
(21) International Application Number: PCT/US99/20449		(74) Agents: FENSTER, Paul et al.; c/o Anthony Castorina, Suite 207, 2001 Jefferson Davis Highway, Arlington, VA 22202 (US).	
(22) International Filing Date: 8 September 1999 (08.09.99)			
(30) Priority Data: 60/100,030 11 September 1998 (11.09.98) US		(81) Designated States: AE, AL, AM, AT, AU, AZ, BA, BB, BG, BR, BY, CA, CH, CN, CR, CU, CZ, DE, DK, DM, EE, ES, FI, GB, GD, GE, GH, GM, HR, HU, ID, IL, IN, IS, JP, KE, KG, KP, KR, KZ, LC, LK, LR, LS, LT, LU, LV, MD, MG, MK, MN, MW, MX, NO, NZ, PL, PT, RO, RU, SD, SE, SG, SI, SK, SL, TJ, TM, TR, TT, UA, UG, US, UZ, VN, YU, ZA, ZW, ARIPO patent (GH, GM, KE, LS, MW, SD, SL, SZ, UG, ZW), Eurasian patent (AM, AZ, BY, KG, KZ, MD, RU, TJ, TM), European patent (AT, BE, CH, CY, DE, DK, ES, FI, FR, GB, GR, IE, IT, LU, MC, NL, PT, SE), OAPI patent (BF, BJ, CF, CG, CI, CM, GA, GN, GW, ML, MR, NE, SN, TD, TG).	
(63) Related by Continuation (CON) or Continuation-in-Part (CIP) to Earlier Application US 60/100,030 (CON) Filed on 11 September 1998 (11.09.98)			
(71) Applicant (for all designated States except US): GENE LOGIC, INC. [US/US]; 708 Quince Orchard Road, Gaithersburg, MD 20878 (US).			
(72) Inventors; and (75) Inventors/Applicants (for US only): STEWARD, Keith, Leroy [CA/US]; 13501 Giant Court, Germantown, MD 20874 (US). SHI, Qin [CN/US]; 10307 Royal Woods Court, Gaithersburg, MD 20879 (US). CARIASO, Michael, Contento [US/US]; 12652 Grey Eagle Court #12, Germantown, MD 20874 (US).		Published Without international search report and to be republished upon receipt of that report.	

(54) Title: GENOMIC KNOWLEDGE DISCOVERY

(57) Abstract

A method of genomic data discovery, comprising: (a) providing a gene data base comprising at least 10 genes; (b) selecting one of said at least 10 genes; (c) discovering knowledge for said selected gene; (d) repeating said (b) and (c) for a plurality of said genes; and (e) repeating said (b)-(d) a plurality of times such that knowledge is discovered substantially in parallel for all the selected genes. Preferably, (b)-(e) are performed substantially without human intervention. Preferably, knowledge discovery utilizes a large number of databases and inference rules to analyze data queried from the databases.



GENOMIC KNOWLEDGE DISCOVERY

FIELD OF THE INVENTION

The present invention relates to automated knowledge discovery and, in particular, to knowledge discovery and data retrieval of genomic related information.

BACKGROUND OF THE INVENTION

Genomic research is currently an expanding field of endeavor. It is expected that by the year 2003, the entire human genome will be mapped out, listing an expected 100,000 different genes. One important customer of this data is the pharmaceutical industry, typically for use as "drug leads". When a new drug is to be developed for a disease, one approach is to discover a gene which interacts with the disease and, then, design a drug which modifies the disease by affecting the gene or proteins generated by the gene.

In earlier days of genomic research, genes were discovered by sequencing expressed mRNA fragments and using these fragments to identify gene locations, using biological methods. Once the gene was located, various methods, usually biological, were used to sequence the complete gene and determine its biological function. The relative amount of available biological data was quite small, so that cross-correlation between data sources was not performed often. These tasks took a considerable amount of time.

In recent years, many methods have been developed for (automatically) generating large amount of biological data, including for example, automatic analyzers, large scale gel-chromatography, high throughput expression profiling and DNA chips. In addition, the amount of research (e.g., micro-biology, neuro biology and structural biology) is also increasing. It is estimated that the volume of data currently doubles every 1.2 to 5 years (depending on the type of data). Thus, current genomic research is characterized by extensive searching through databases for correlations with available data. In parallel, other steps of drug discovery have been automated, for example, automatic screening of drug candidates.

A current paradigm for drug research is as follows:

- (a) identifying genes which mediate and/or are involved in diseases and/or biological processes of interest;
- (b) establish identity and/or gene class of the identified genes;
- (c) selecting genes with the most useful biological properties;
- (d) screening for a lead compound which modulates the selected gene(s) or gene product(s);
- (e) selecting related compounds which are even better; and

situation matches a certain rule, that rule fires, the result of which is the performance of the action.

International Business Machines Corporation, of Armonk, N.Y., is developing "Agent Builder" a commercial system for creating intelligent agents, in which each agent includes a reasoning engine based on inference rules.

SUMMARY OF THE INVENTION

The invention described herein, in some embodiments thereof, allows a fundamental change in the way genes are evaluated for the purposes of drug target discovery. Currently the dominate paradigm for gene evaluation requires that research analysts use one database or tool at a time, on one gene at a time and integrate and collate the resulting information before moving on to the next database, tool, and/or gene. Although there are extensive numbers of databases and tools at the disposal of biologists today, only a few are typically used because of the effort involved. After several such actions, the integrated data are correlated in order to detect implications in the data, and in order to prioritize the genes for further possible drug target consideration. Given the current explosion in genomics data (primarily genes), brought about by the emergence of high throughput genomics technologies, together with the requirement for exhaustive searching and analysis of the data before the very rare drug targets can be found, new methods for integrating the data presently available are probably required to assure the success of genomics-based drug discovery.

In accordance with a preferred embodiment of the invention, success in this genomics-driven drug development field is enabled by a break-through in the way this vast amount of complex data are accessed, integrated, and correlated to discover actionable conclusions (knowledge). Unfortunately the most popular 'data mining' and other tools available in other problem domains are often not useful in the genomics context, because of the scale, complexity, and heterogeneity of the data. The present invention, in some preferred embodiments thereof, provides a novel and powerful solution to the problem in the form of an 'intelligent' agent architecture which is specifically designed and optimized for the nature of the problem and the associated data.

One object of some preferred embodiments of the invention is to provide an automated system which enables a biological researcher to deal with challenges caused by problematic data. Preferably, the system aids the researcher to overcome challenges caused by one or more of the following: scale, updating, complexity, heterogeneity and garbage. Preferably, the system relieves the researcher of a need to become an expert in manipulating the bewildering range of available data sources and tools.

(i) analyzing the integrated information and/or comparing it to other information, yielding inferences and conclusions, i.e. knowledge;

(j) facilitating conclusions and rational actions from the knowledge; and

(k) iteratively repeating the cycle (from step b). In some embodiments, at least some of the above steps may be varied, skipped and/or be performed in a changed order. In a preferred embodiment of the invention, the decision to skip or change an order of a step is made by a knowledge discovery system, preferably depending on circumstances as recognized by the system.

One aspect of some preferred embodiments of the invention is that the cycle is performed on one gene token at a time. A complete round of data mining and knowledge discovery cycles preferably includes applying the cycle once for each gene token in the gene token database. Thus, the state of knowledge for all the gene tokens in the gene token database is advanced substantially in parallel. Typically, one does not know at the start of the process which gene will yield a significant amount of knowledge.

An aspect of some embodiments of the present invention relates to goal setting in a knowledge discovery system. In a preferred embodiment of the invention, the system does not work towards a single predefined goal. Rather, a "target area" of desired results is defined and any hit within the target area is considered a success. Alternatively or additionally, any knowledge accumulated during the discovery process may be useful, for example for advancing the state of knowledge about particular genes. The target area is preferably not explicitly defined. Rather, useful properties of results falling in the target area are defined. Alternatively or additionally, the system encodes heuristic rules (preferably as inference rules) for advancing knowledge towards the desired target area.

An aspect of some preferred embodiments of the invention is that each gene token is associated with one or more schemes (or frames, slots or other similar AI constructs). Thus, missing and/or possibly available information can be determined from the scheme. In a preferred embodiment of the invention, flexible and/or extensible schemes/frames are used. Thus, widely varying attributes of gene tokens may be accommodated. Another advantage of some embodiments is that different gene tokens will have different extents of knowledge, and the schemes/frames can offer the necessary flexibility to accommodate this. Another advantage of some embodiments is that an automated scheme can be used to represent a significantly more complicated data structure that can be conveniently utilized by a human researcher, for example for setting data requirement goals and/or for analyzing the resulting data. Although a user may be able to enumerate all the possible slots and/or interactions between slots in a

stopping work on a gene token or reporting the gene token, may be triggered when the volume and/or quality and/or nature of the data for a particular gene token reach a particular state.

An aspect of some preferred embodiments of the invention relates to corrective action performed when an error is detected. Data errors may be detected by comparing data from
5 different databases. Alternatively or additionally, the data in a database may not be consistent. Alternatively or additionally, the data in a database may be updated. In a preferred embodiment of the invention, when a knowledge discovery process is complete (or reaches a certain stage) the conclusions are compared to the data used to reach the conclusion. Some conclusions may trigger re-examination of the data, or correlation with other data, in order to
10 detect inconsistencies or errors in the data. Although some inconsistencies are caused by erroneous data, some may be caused by a misinterpretation of available data, too broad generalizations, application of erroneous rules and/or erroneous application of rules. In a preferred embodiment of the invention, when an error is detected in a database, inferences based on the erroneous data may be deleted and/or reapplied. Alternatively or additionally, the
15 data in the erroneous database is corrected, preferably with an indication in the database of the correction. In an external database, the system may generate a communication, such as an e-mail explaining the detected error and reasoning behind the error. When the data source is a programming tool, such an indication may be sent to a programmer and/or a maintainer of the program. When the data source is a laboratory, the error may be sent to laboratory personnel.
20 In a preferred embodiment of the invention, the knowledge discovery system includes a plurality of rules which attempt to explain possible sources of data errors, based for example, on other available data.

An aspect of some preferred embodiments of the invention relates to automatic modification of system behavior, responsive, for example, to self-monitoring results. In one
25 preferred embodiment of the invention, the system selects databases which have a lower determined error rate and/or higher availability for the searched for data. Alternatively or additionally, the system optimizes various parameters thereof, in response to the self monitoring. Alternatively or additionally, the system modifies rule results and/or rule activation, based on the self monitoring.

30 An aspect of some preferred embodiments of the invention relates to prioritization and/or resource allocation. In a preferred embodiment of the invention, various gene tokens are attributed a higher priority than other gene tokens. Knowledge discovery may be performed more rapidly, more often and/or utilizing more computationally expensive tools and/or resources, on tokens having higher priorities. In one example, human resources may be

In a preferred embodiment of the invention, the system automates workflow, for example, by generating work orders and/or prioritizing and/or allocating physical resources, such as laboratories, personnel and/or computer equipment. In a preferred embodiment of the invention, the system can directly or indirectly (through control programs) operate laboratory equipment, for data generation.

An aspect of some preferred embodiments of the invention relates to controlling data generation devices. In a preferred embodiment of the invention, the system may allow lower quality or reliability biological data to be provided. If later gleaned knowledge requires it, higher quality data may be requested by the system, for example via a work request sent to laboratory personnel. Alternatively or additionally, when requesting data the system defines a required quality level. Thus, the system can optimize throughput of laboratories, so that they operate at peak performance, generating data having only as high a quality or reliability as required.

An aspect of some preferred embodiments of the invention relates to automation level. Various levels of automation may be applied in different embodiments of the invention, at one end of an automation scale is a system which, once programmed, does not require any human input. On the other end of the scale is a system which will not operate unless it is continuously being controlled by a human. In between, are systems which report some or all activities, systems which request an OK for some or all activities, systems which notify a user of some or all happenings and/or systems which operate autonomously for a certain period of time, before requiring a user OK to continue. In a preferred embodiment of the invention, the system includes a plurality of automation levels. Preferably, different activities and/or results are handled differently. In one example, limits on resource expenditures may be defined. Going over the limits may require a user intervention. In another example, some results are determined to be important enough to be reported to a user immediately and some less urgent results are accumulated for a periodic report.

An aspect of some preferred embodiments of the invention relates to automatic notification of users, regarding newly available content (e.g., data, information and/or knowledge). An important consideration is not to overload a recipient with too much information. Additionally, information and/or knowledge should preferably be presented only when it's content and/or confidence level meet certain criteria. Additionally it should be noted that there is often no "requester" for information. Rather, the information and/or knowledge may have simply come up. In a preferred embodiment of the invention, a set of inference rules is defined to determine if data, information and/or knowledge should be present and who are

example, the prioritization of the token is modified responsive to an analysis of the biological data.

An aspect of some preferred embodiments of the invention relates to various methods of prioritization of gene tokens for usefulness in a particular application. Preferably, the tokens
5 are ranked according based on biological data. Alternatively or additionally, other data pertaining to the application may be utilized. Preferably, the application for which the genes are ranked is a drug discovery application. One method is based on the biological relevance of gene tokens, for example, matching across biological spatial/temporal/event/genomic map dimensions, described above. Another method is based on pharmaceutical tractability. Some
10 types of proteins may be more difficult to design drugs for; for example due to functional and/or structural characteristics and/or due to a difficulty in properly localizing the drug. Alternatively or additionally, some proteins may be similar to housekeeping proteins and/or significant other proteins, increasing the possibility of side effects. Alternatively or additionally, some proteins may be less suitable targets due to their expression profile and/or
15 existence of parallel pathways, which operate even if the protein is knocked out. Another method is based on experimental tractability. As can be expected, experimental validation of a drug target is simpler if the gene has a close homologue in a different animals (for example a mouse). Also, proteins which cross the cell membrane may be easier to work with than proteins which remain in the cell nucleus. In a preferred embodiment of the invention, the
20 above considerations are encoded as inference rules, background assertions, and/or criteria.

An aspect of some preferred embodiments of the invention is the utilization of knowledge discovery cycles including a plurality of points at which inferences may be applied. Generally, different inference rules are applied at each point. Preferably, the points at which
inference rules may be applied include one or more of: data retrieval goal definition, data
25 source selection, search results analysis and extraction, data integration, reporting/workflow to users and prioritization of gene tokens.

An aspect of some preferred embodiments of the invention is that inferences are preferably based on a meaning of terms (semantics), rather than on an exact word match. In a preferred embodiment of the invention, terms are interpreted using a semantic network, whose
30 content is preferably derived, at least initially, from UMLS (unified medical language system). Thus, the terms IL-2 and interleukin 2 will be interpreted as the same concept. Alternatively or additionally, the broadening of a term may be applied at a query construction stage ('query expansion'). Preferably, the broadening for a particular database is limited to exclude terms

preferably enables an adapter to be rapidly reconfigured to begin working with an internal copy of a database that was once accessed externally, for example over the Internet and/or a dial-up connection.

5 An aspect of some preferred embodiments of the invention is that an adapter registers itself, when it becomes available to the agent. Thus, inference rules for selecting a relevant data source have available a list of available adapters, access properties (such as time and cost), field list, range of available data and/or other information which may aid in selecting a data source for retrieving data.

10 An aspect of some preferred embodiments of the invention relates to communication between adapters and an agent. Preferably, the adapters each communicate with the inference part of the agent using their own mailbox. Alternatively or additionally, a different mailbox is provided for each type and/or priority of message. Alternatively or additionally, a single mailbox is shared among some or all the adapters. In a preferred embodiment of the invention, each adapter operates as a separate computational thread, thus, the agent is not required to wait
15 for requested information to be retrieved. Preferably, the adapters can communicate between themselves, for example, to coordinate usage of communication resources. Alternatively or additionally, such communication is performed through the inference engine.

In some cases, the agent may not wait for any information. Thus, some cycles might be devoid of retrieving results and contain only information analysis. Other cycles may include
20 only data requesting and no analysis, as no new data is available. Preferably, the agent does pause for data retrieval at least a limited amount of time. Alternatively or additionally, the agent pauses responsive to a number, content and/or time stamp of outstanding data requests. Alternatively or additionally, an agent may pause until a particular data element is retrieved. Preferably, such pausing is dictated by suitable inference rules. Alternatively or additionally,
25 an agent may select a faster- or a slower- accessed version of database, for example if a same database is accessible both on a local drive and externally. Such a selection may be responsive to an update status of the database and the instant importance of receiving up-to-date data. Alternatively or additionally, such a selection may be responsive to communication bottle necks. In one example, the inference engine may provide an adapter with a set of data sources,
30 data from any of which may be sufficient. In another example, an agent may request data from a lower quality source, to be received faster and, in parallel, request data from a higher quality source (for example an on-line mirror of a local database), to be used to possibly correct inferences, when the data arrives.

In a preferred embodiment of the invention, said at least 10 genes comprises at least 50 genes. Alternatively or additionally, said at least 10 genes comprises at least 100 genes. Alternatively or additionally, said at least 10 genes comprises at least 300 genes. Alternatively or additionally, said plurality of genes comprises at least 20% of said at least 10 genes. Alternatively or additionally, said plurality of genes comprises at least 50% of said at least 10 genes. Alternatively or additionally, said plurality of genes comprises at least 80% of said at least 10 genes.

There is also provided in accordance with a preferred embodiment of the invention, a method of genomic knowledge discovery, comprising:

- 10 determining at least one required data element for at least one gene;
- querying a plurality of at least 50 databases for said at least one required data element;
- receiving responses from said databases; and
- analyzing said responses to increase knowledge for said at least one gene.

Preferably, said at least 50 databases comprise at least 100 databases. Alternatively or additionally, said at least 50 databases comprise at least 300 databases. Alternatively or additionally, said databases are all queried for a same data value. Alternatively or additionally, said method is performed automatically.

There is also provided in accordance with a preferred embodiment of the invention, a method of automated knowledge discovery, comprising:

- 20 continuously operating a knowledge discovery cycle comprising:
 - querying a database to receive data; and
 - performing inferences on said data to generate knowledge; and
 - re-evaluating said inferences when said database is modified

Preferably, said cycle is continuously operated over one week. Alternatively or additionally, said cycle is continuously operated over one month. Alternatively or additionally, said cycle is continuously operated over six months.

There is also provided in accordance with a preferred embodiment of the invention, a method of genomic knowledge discovery, comprising:

- 30 (a) selecting a gene token;
- (b) determining data requirements for said gene token;
- (c) requesting and receiving data responsive to said data requirements;
- (d) analyzing said received information to increase knowledge for said gene token; and
- (e) repeating (b)-(d) for the same gene token, at least 50 times.

Alternatively or additionally, said application comprises drug discovery. Preferably, said application specific rules comprise biological relevance rules. Preferably, said biological relevance rules comprise rules which match said gene tokens to a disease model across a plurality of biological dimensions.

5 Alternatively or additionally, said application specific rules comprise rules which determine experimental tractability. Alternatively or additionally, said application specific rules comprise rules which determine pharmaceutical tractability. Preferably, pharmaceutical tractability determination comprises an indication of ease of finding an effective pharmaceutical. Preferably, pharmaceutical tractability determination comprises an indication
10 of ease of finding a pharmaceutical with a low level of side effects.

There is also provided in accordance with a preferred embodiment of the invention, a method of genomic information analysis, comprising:

providing a first model of a biological relationship which interrelates a first plurality of genes or proteins;
15 providing a second model of a biological relationship which interrelates a second plurality of genes or proteins; and
applying inference rules to said first and second models to infer missing information.

Preferably, said applying comprises determining data needs for at least one of said models, based on said applied inference rules. Alternatively or additionally, said biological
20 relationships comprise enzymatic pathways. Alternatively or additionally, said biological relationships are in different species.

There is also provided in accordance with a preferred embodiment of the invention, a method of automated genomic knowledge discovery, comprising:

analyzing a gene token to determine required data;
25 selecting at least one suitable human expert; and
querying the at least one selected human expert for the required data.

There is also provided in accordance with a preferred embodiment of the invention, a method of automated genomic knowledge discovery, comprising:

analyzing a gene token to determine required data; and
30 automatically generating a work order to a laboratory to generate the required data.

BRIEF DESCRIPTION OF THE DRAWINGS

The present invention will be more clearly understood from the following detailed description of the preferred embodiments of the invention and from the attached drawings, in which:

In a preferred embodiment of the invention, a gene token may have associated with it an alert request, so that if data becomes available, it is forwarded to the gene token. Automatic notification is a feature provided in many databases and is usually based on a query which is re-run periodically at the database. If any new data becomes available, the user of the database (in this case preferably the knowledge discovery system), is notified. In one example, a gene token may "set up" an alert on sequence information for the gene. Alternatively or additionally to such an alert function being performed by outside data sources, the knowledge discovery system can also be used to periodically search relevant data sources and automatically notify users (or hibernating tokens, on which no action is to be performed pending new information) of newly available data.

In step 26, data sources which can respond to the queries are identified. In a preferred embodiment of the invention, each available data source is registered in a central registry. Preferably, the same data is queried from a plurality of data sources, to overcome problems caused by erroneous data, missing data and/or slow response times.

In step 28, the selected data sources are queried. In a preferred embodiment of the invention, the queries are adapted to match the particular conventions and/or formats of the selected data sources. In a preferred embodiment of the invention, the queries are adapted by specifying the data needs in terms of canonical concepts and translating them into terms and a syntax which are understood by the data source. This syntactic and semantic translation of the data needs (data mining goals) into queries is preferably carried out by each data source adapter. Alternatively, if the adapters are hierarchically organized, at least some of the translation may be performed a single time for a set of related adapters.

In step 30, the results are retrieved from the data sources. In a preferred embodiment of the invention, the results are parsed in order to determine their content. Some data sources return results in a relational format. Others however, return strings, for example "gene RAS was found in 20% of tissue type X". Such a string is preferably parsed and the term "found in ... tissue type" interpreted to mean the same as $\text{Occurrence}(\text{RAS}, \text{X})=20\%$ and "20%" (where the query was "in what percentage of tissue X is RAS found"). In addition, some formats are standard in the field of genomics, for example pathways are usually written as "XXX pathway". In a preferred embodiment of the invention, if the results are ambiguous, semantic mapping techniques, described below are used. Alternatively or additionally, the content of a result may be used to constrain a description of a result. For example, a field value which is formatted similar to a radiation hybrid map location may be used to interpret a field name "map location" as "radiation hybrid map location".

inference rule set 58. Inference rules are described in greater detail in a separate section, below.

In a preferred embodiment of the invention, database 40 itself is also one of the databases queried in step 28 of Fig. 1, since two gene tokens may refer to similar or overlapping things.

In the following example, a candidate GL375 is the seed. An example of data which is in the scheme is an association with a sequence EST 9224. Querying for a sequence matching the EST might yield that there is a match with osteoblast-specific GPCR and a map location of the gene. A second round of queries (at a later time) might yield the fact that Osteoporosis has been genetically linked with that same map location. An example of knowledge generation is that GL375 is a good candidate to account for osteoporosis.

In Fig. 1, not all the steps are necessary in order to increase the level of knowledge for each gene token. Additionally or alternatively, the order of steps may be modified. In one example, identification of suitable data sources is not performed. As a result, many databases may return empty responses. In another example, the information from the data sources is not integrated. Instead, the first arriving data or the data with the highest reliability value (based on the database identification and/or provided by the database) are used. In another example, data needs may be defined after the information is analyzed, so that a token is always associated with outstanding data needs.

Fig. 3 is a flowchart of a process of updating a gene token database in accordance with a preferred embodiment of the invention. Preferably, the knowledge discovery cycle of Figs. 1 and 2 are continually being performed. In Fig. 3, a round comprises performing one knowledge-increasing step for each gene token in database 40. Thus, the level of knowledge for the entire set of gene tokens is always advancing. As can be appreciated, it is not usually known in advance which of the gene tokens will yield a viable drug target.

In a preferred embodiment of the invention, the rounds are driven by selection of gene tokens from database 40. Alternatively or additionally, the rounds may be driven by external events, for example, an updating of data or a specific user request.

Fig. 4 is a schematic illustration of an integration of a knowledge discovery system 70 in an industrial setting, in accordance with a preferred embodiment of the invention. One preferred industrial setting is gene discovery and drug lead generation, in which genes are selected based on their suitability to become drug targets.

As shown in the Fig., system 70 may be connected to management 70, one or more external data sources 74, one or more internal data sources 75, one or more laboratories 76, a

and/or taking into account special requests. Functional characteristics include, for example, a managerial level, a person being a head of a project (responsibility level), secrecy considerations (compartmentalization), and/or a person being part of a project, part of a research group or being a more generally open fielded individual. Thus, an interesting
5 discovery regarding an osteoporosis associated gene may be communicated to a project working on drug leads for osteoporosis, to a researcher who generated a significant item of the information leading to the discovery, to a researcher in the field of body fluid calcium levels, to a manager (based on the level of interest of the item) and/or to a worker who requested to be kept up to date on osteoporosis related information. Personal characteristics include, for
10 example, an annoyance level, number of items of information to be sent per week, a known workload, a personal interest field and/or interest level threshold.

In a preferred embodiment of the invention, information is sent by e-mail preferably including hypertext/HTML email with links to richer or interactive communications. Alternatively or additionally, the information is published on a network, for example using an
15 Intranet or an Internet. Alternatively or additionally, the information is sent by fax. Alternatively or additionally, the information is sent by voice, for example directly into voice mailboxes of recipient.

In one embodiment of the invention, system 70 is formed of two sub-systems, a knowledge management system and a knowledge discovery system which is treated by the
20 knowledge management system as one of its users.

GENERAL SYSTEM STRUCTURE

Fig. 5 is a schematic block diagram of a knowledge discovery system 70, useful for Fig. 4, in accordance with a preferred embodiment of the invention. the functionality of system 70 is preferably provided by one or more agent instances 102 which perform the above
25 described knowledge discovery process. Each agent preferably comprises an inference engine 104, which is used to manipulate knowledge representations and/or assertions pertaining to gene tokens and/or diseases, and one or more adapter 110, which are used to transfer data to and from the agent and/or to otherwise communicate with an external world, for example, to control one or more tools 116. The collection of adapters are preferably contained in an
30 interface layer 109, which preferably comprises a program interpreter (preferably Perl), which is optimized for text processing, inter-process communication and/or rapid development. Inference engine 104 is preferably written in CLIPS. An agent backbone 106 is preferably provided to link inference engine 104 with adapters 110 and to provide an operational setting for agent 102. Preferably, backbone 106 includes a mailbox 108, for adapters 110 informing

- (e) terms used in output (e.g. "is a variant of");
- (f) terms useful in queries;
- (g) formatting of particular types of values;
- (h) ranges of values in fields and/or coverage of a database;
- 5 (i) a confidence level (per database or per different parts of a single database);
- (j) expected response time; and/or
- (k) update status.

In a preferred embodiment of the invention, data source registry 112 is embodied, at least in part, as assertions. Alternatively, data source registry 112 is provided by a mediator
10 software component, possibly an adapter, which provides a list of relevant databases, when queried. In a preferred embodiment of the invention, the information in data source registry 112 also affects query formulation, for example by an adapter providing rules or assertions which indicate what data could be requested. Alternatively or additionally, such inference rules may be used to expand a scheme associated with a gene token. For example, a gene token
15 scheme may include only sequence information. If an adapter is provided which matches sequences with protein acidity sensitivity, a new slot entitled "protein acidity sensitivity" may be added to the scheme. Such an event is preferably also indicated to programmer 86 (Fig 4), so that appropriate information analysis rules may be added.

It should be appreciated that a suitable adapter can communicate with any type of
20 information source, including data sources, e-mail and/or data manipulation tools. In one example, an adapter may be used to run a BLAST homology test. In another example, a different adapter may be used to perform some or all of the information analysis step. Alternatively, an adapter may be used to execute a program which performs any of the steps described with reference to Fig. 1. In another example, a special adapter may be used to match
25 schemes to existing gene tokens. Thus an adapter may send and/or receive various types of data, including data objects of agent 102, such as schemes or rules.

In a preferred embodiment of the invention, adapters are written in a text-processing language. Preferably the language is interpreted. Preferably, the language is Perl. In a preferred embodiment of the invention, adapter design takes into account similarities between data
30 sources. In one example, the adapters are designed in a modular manner, so that only those modules which are different need to be programmed. Additionally, such modular design aids in updating adapters when data source formats change. Additionally, such modular design aids when a database is to be accessed in a plurality of methods, for example over the Internet and from a local CDROM drive.

parsing the results of a query, so that information contained in the database's records are returned to the inference engine in a canonical or standardized vocabulary. The semantic mapping may be performed by the adapter for the data source. Alternatively or additionally, the semantic mapping is performed by the inference engine.

5 In a preferred embodiment of the invention, the semantic mapping is performed using a comprehensive database of biomedical terms 114, for example, initially populated with content from the Unified Medical Language System (UMLS) knowledge base, available from the National Library of Medicine. Preferably, the database is converted into inference rules and/or assertions. Preferably database 114 is continuously updated, for example by users or by system
10 70 itself. In one example, system 70 updates semantic mappings based on a query result in which such mappings are explicitly stated, for example in a text based database: "...cardiac muscle (e.g., papillary muscle)...", may be interpreted by system 70 to mean that papillary muscle is a type of cardiac muscle, at least in that database.

In some cases, a syntactic analysis may aid in the semantic analysis. For example as
15 described herein with respect to radiation hybrid map locations, a format used in a field may be used to guess at or narrow the range of possibility to a semantic meaning of a field name. For example, a field that is labeled 'Map Location' (of which there are several possible types) and containing a value of '21cR' may be mapped to the concept of 'rad-hybrid-map-location = 21'.

20 In a preferred embodiment of the invention, system 70 is implemented on a LINUX machine with at least 64Mb. Preferably, a plurality of agents are simultaneously executed on a single machine. Alternatively or additionally, a plurality of agents are executed on a plurality of networked machines. In a preferred embodiment of the invention, the machine and/or the network are selected with redundancy in mind, so that system 70 will be able to continue
25 normal operation even if a hardware or software component fails. One advantage of LINUX computational platforms is their excellent performance/price ratio. Thus, large numbers of agent instances may be effectively deployed to analyze a large number of gene tokens.

In a preferred embodiment of the invention, each agent is executed as a multi-threaded application. Preferably, each adapter is in a separate execution thread, as are the backbone,
30 interface layer, and/or the inference engine. Hence, if one adapter is stuck, the rest of agent 102 is not significantly affected, leading to significant robustness. Alternatively or additionally, inference engine 104 may also be multi-threaded, especially to support multiple simultaneous inference chains. Also, it should be noted that an inference rule may depend on the execution of an adapter. Alternatively or additionally, the agent may be implemented in a multi-process

phasic distribution. Conversely, two gene tokens may be merged, if there is a high enough match between them. Alternatively or additionally, related genes may be clustered and/or otherwise associated, in that they receive similar treatment.

In a preferred embodiment of the invention, schemes are created by users. Alternatively
5 or additionally, schemes may be updated automatically, for example in response to a query result in which data fields, which are not supported by the scheme, are retrieved or as a result of parsing a free-form query response, which response contains a relationship not covered by the existing scheme. The association of one or more particular schemes with a gene may be performed manually. Alternatively or additionally, system 70 may suggest and/or associate the
10 schemes. Alternatively or additionally, system 70 may modify the schemes and/or modify the scheme associations, for example based on inference rules which analyze an existing state of knowledge about a gene token.

In a preferred embodiment of the invention, one token may message a second token, for example if a piece of data can only belong to one of the tokens. The result of such a
15 message may be an indication in one token that a previously correct data is not considered to be suspect. Another possible result is the activation of a mediation process, in which differences between the two tokens are settled. Another type of message is when the tokens belong to the same family of tokens and one token comes up with information which may be of interest to a second token. Preferably, the tokens in a family are explicitly linked, for
20 example as part of their scheme, so that not all actions at one token will trigger an action at a related token. The above messages are preferably implemented by setting attributes of schemes associated with gene tokens, which attributes may initiate the firing of "message handling" rules.

It should be noted that in system 70 it is expected that one gene token be "researched"
25 based on knowledge advances in a second token. For example, an osteoporosis token may utilize knowledge gleaned while researching the maintenance of serum calcium level. Such chaining of inferences between two, three or more tokens may be performed automatically, as described herein, for example by treating the gene token data base as a data source for data retrieval. In addition to chaining, other "collaboration" configurations are expected, for
30 example, two or more gene tokens co-advancing, each one building on the other; star configurations where a single gene token is fed from a plurality of gene tokens; token rings, where each token is dependent on a previous token; and more complex configurations, for example as can be described using a directed graph or tree. In a preferred embodiment of the invention, such interactions are detected by system 70 (preferably using suitable inference

In a preferred embodiment of the invention, system 70 includes several possible rule sets and background assertions (knowledge). When a user sets up a task, the user preferably selects one or more rule sets to use in the make up of an agent 102. In a preferred embodiment of the invention, some rule sets may be specialized per task or per scientific outlook. 5 Alternatively or additionally, a rule set may be crafted and/or modified so that it is particular to a disease and/or a disease model. Alternatively or additionally, a rule set may include inference rules which fire on schemes matches. Alternatively or additionally, the scheme may be embodied, at least in part, as rules.

In a preferred embodiment of the invention, one of the first rule types applied to a gene token is used identifying the gene token. In one example, a gene token is provided as a EST. 10 The applied rule preferably identifies that a complete sequence is missing and searches in relevant data bases for sequences with homologies to the EST. The series of data mining and knowledge discovery goals may follow a particular line of reasoning, or a particular hypothesis exploration, which may possibly be a characteristic of the agent, gene token, associated 15 schemes and/or rule set.

In a preferred embodiment of the invention, inference rules are used to set priorities between gene tokens. As explained herein, priorities may be used, inter alia, for resource allocation and/or for selecting recipients for data. Some gene tokens may be marked as not being active or may be advanced only once every few rounds. Preferably, different resources 20 may have different priorities associated with them, per gene token. For example, two gene tokens may have opposing time priority and money priority.

In a preferred embodiment of the invention, prioritization is based on a score. The determination may be made by comparing the score to a threshold or by comparing the score to scores of other gene tokens. In a preferred embodiment of the invention, an effective score 25 also takes into account changes in score as a function of expenditure of time, money, cycles and/or other resources. Possibly, the ranking score for a particular gene token is based on multiple components. Preferably, gene tokens whose score does not increase, are provided with a reduced priority. Alternatively or additionally, gene tokens are removed from consideration if knowledge discovery may be terminated, for example, if they have reached a 30 complete description, have been identified as useful candidates, have been identified as house-keeping genes and/or if no further data has been found.

In a preferred embodiment of the invention, a score is generated based on matching to ranking rules. For convenience, these rules may be divided into the following types:

As an example of research tractability, some genes families have been more heavily researched. Thus, more relevant biological information may be expected to be available. The volume of possibly relevant information is often an important consideration in knowledge discovery.

5 As an example of financial tractability, the importance of developing a pharmaceutical is often a function of the market for the drug. If only a small market is available, the payoff may be too small to take a risk. The risk is preferably assessed using the other prioritization rules and/or by providing marketing information.

10 It should be noted that the above rules may also be used to select research areas in which human researchers have not mined out, since they were expected to be unproductive.

In a preferred embodiment of the invention, inference rules are used to draw analogies between biological situations, in order to point at missing data and/or offer explanations. In one example, two pathways may be compared. For evolutionary reasons, similar pathways often utilize homologous genes. Thus, if one pathway is known, it is possible to expect to find
15 homologous genes in a similar pathway (in the same animal or even between species). This type of homology may also be used for a negative purpose, for example to determine an increased risk of side effects.

AUTOMATION LEVEL

System 70 may be operated at various levels of automation. Alternatively or
20 additionally, different components of system 70 may be operated at different levels of automation. In general, automation levels vary between a completely manual system, when a person is required to perform all the activities and make all the decisions and a completely automatic system, where the system does as it will and does not even notify a person after the fact. In a preferred embodiment of the invention, system 70 operates at an intermediate level.
25 For example, the system may require human intervention for some decision, perhaps only the final decision after much inferencing and decision making has already been performed by the system. Also, the system may notify a person of a made decision. Also, a person may affect the operation of the system directly, for example by command, or indirectly, for example by modifying a control-type inference rule or by setting a goal.

30 In various preferred embodiments of the invention, any of the above described activities may be performed manually, automatically or semi automatically, meaning that system 70 aids a person in performing the step. Generally, it is preferred to limit human intervention, at least to those activities which i) are too difficult to program a computer for, or ii) for an ultimate decision involving the commitment of large amounts of human and/or other

USER INTERFACE

One important aspect of user interface has already been mentioned, matching data sent to a user to the user's desire, ability and/or job function (responsibility). In a preferred embodiment of the invention, system 70 includes a user model, to better guess if, when and how content should be sent to a user. Preferably the user model is based on psychological characteristics of a user, possibly entered by a user himself. Alternatively or additionally, the user model is based on a task being performed by the user and/or on his responsibilities. Thus, new discoveries may be reported immediately to a researcher, while such new discoveries will be ranked and reported to a manager once a week. Additionally or alternatively, the model for a particular user could be learned during interactions with the user over time.

In a preferred embodiment of the invention, one or more of the following human interfaces are provided:

- (a) e-mail, sending and receiving;
- (b) WWW pages, for filling forms and/or for receiving and/or transmitting graphics-rich and/or interactive responses;
- (c) command line prompt, for commands;
- (d) menus, preferably for commonly used commands; and/or
- (e) by receiving and/or generating electronic files.

In a preferred embodiment of the invention, the type of interaction can be immediate or less so, preferably depending on an urgency; for example:

- (a) setting up tasks to be performed, for example gene tokens and/or inference rule sets, or changing inference rules, which actions do not generate a response in the near future;
- (b) setting up breakpoints and/or report-points;
- (c) requesting a status report and or a display of data base contents;
- (d) scripts to be executed under certain circumstances;
- (e) other ad hoc commands, for example requesting a roll-back of an inference chain; and/or
- (f) real-time or semi real-time question/answer type interactions.

In a preferred embodiment of the invention, results are displayed as a ranked list. Preferably, a graphical interface is provided for a user to examine the results, preferably using a point-and-click interface. Such examination may include displaying an inference chain, displaying activated rules, displaying confidence levels, displaying raw data, displaying data sources and/or displaying inconsistencies, within the gene token and/or with other sources of data. Preferably, the data is displayed in a hierarchical manner. for example, clicking on a gene

many parts of an organization. In a preferred embodiment of the invention, some of the activities of system 70 are directed towards optimizing the contact and/or providing feedback to the parts of the organization. In one example, system 70 tracks (for example by using "reader" response messages) the effectiveness and/or suitability of communications to users.

5 Such feedback, for example can be used to adjust the 'annoyance level' in a user's user model, as described above.

In another example, system 70 may perform QA/QC (quality assurance/quality control) functions. Typically, a significant amount of data generated by the organization may be used to feed system 70. Any problems with the data are preferably detected by system 70.

10 Alternatively or additionally, system 70 includes inference rules for explaining inconsistency by indicating a problem with data. Also, system 70 can cross-correlate internal results with data from external databases. Also, system 70 can compare old conclusions against new data. The notification of appropriate personnel of these problems by system 70 has already been described above.

15 In a preferred embodiment of the invention, an activity of system 70 is resource allocation. In some cases, the only "noticeable" resource used by agents 102 is computer time. In addition however, system 70 may have other resources, for example money, for paying for data from fee-based data bases, work-hours and laboratory-hours, for performing work by people and/or machinery, communication bandwidth, disk storage space and time, for example

20 deadlines or time windows where certain activities may be performed. Another type of resource is a database which can only be accessed one-at-a time.

In a preferred embodiment of the invention, resource allocation is reasoned over via the inference engine, and these resources are allocated based on priority between the gene tokens. Additionally or alternatively, each gene token may have associated therewith an upper spending

25 limit and/or an upper spending limit per cycle. In a preferred embodiment of the invention, tokens are evaluated using a round-robin mechanism. Preferably however, a timer is set so that the data retrieval and knowledge discovery cycle is kept below a maximum duration. Alternatively or additionally, other resource allocation mechanisms may be used. For example, if work on a particular agent is hanging due to input being awaited, control may pass to a

30 different token or a different agent on the same computer.

In a preferred embodiment of the invention, configuration control extends to data, information knowledge and/or rules. In a preferred embodiment of the invention, each gene token has stored therewith (perhaps in the gene token database) a snapshot of the current data mining / knowledge discovery state, allowing the agent to return to this state at a later time.

- (b) relative access to different databases;
- (c) error rate in different databases;
- (d) acceptance rate of inferences by users;
- (e) data sources which do not respond (or a temporal profile of response times); and/or
- 5 (f) positive and/or negative contributions of data sources and/or of individuals to the success of knowledge discovery.

Preferably, the results of the self-monitoring are presented to an operator. Alternatively or additionally, some of these results may be used by the system to optimize its performance, for example by using erroneous databases less often and/or by using databases which had been
10 previously neglected.

DEALING WITH ERRORS

In a preferred embodiment of the invention, any action or inferences taken by system 70 is logged, so that the source of errors (and good results) may be traced. Not only might this allow the system to analyze its past performance and induce new rules based on this analysis,
15 but also preferably, when a user checks a gene token and/or a current state of a system, the user also checks the data used to reach the inferences – this is possibly in addition to checking the inference rules themselves. Often, a piece of data which is accepted by system 70 may be immediately identifiable as erroneous by a user. Alternatively or additionally, the user analyses the confidence level of the system in data, information and/or knowledge. Alternatively or
20 additionally, a user can selectively view other portions of system 70, for example those which formulate queries, set data requirement goals and/or parse. If these portions are implemented as inference rules, a user can preferably selectively view only those inference rules. If these portions are not implemented as inference rules a user can preferably view source code, parameter settings and/or configuration files for these portions.

25 In a preferred embodiment of the invention, system 70 includes automatic error correction mechanisms. Preferably, when an error or an inconsistency is detected, system 70 attempts to determine where the error occurs, for example in which data source. Alternatively or additionally, system 70 isolates the error and freezes inferences involving the error, for example, at least until a user intervenes. In one example, if three data sources return sequence
30 information, which do not match, a vote may be cast to select the correct sequence data. As a result, the contig and/or full-length sequence is preferably identified. Further inferences may yield a protein, protein family, protein homology and/or protein function.

In some cases, correctness may be determined by allowing the system to continue using the possibly erroneous data and then attempting to determine, at a later time, if indeed

In a preferred embodiment of the invention, system 70 may perform testing activities, to compare two variants of a rule, parameter setting, data interpretation, threshold value and/or other modifiable items in system 70. Preferably, system 70 tests a hypothesis, regarding which variant is to be preferred, by executing two or more variants, preferably simultaneously. The
5 variant which achieved the better results or the variant which did not cause any (or less) inconsistency in stored content is preferably considered to be the better variant.

SOME POSSIBLE MODIFICATIONS

Although many functions of system 70 may be achieved using inference rules, in some preferred embodiments of the invention, these functions may be achieved using other means.
10 Such means include scripts, preferably in an interpreted language; production rules with associated databases of situations; and/or procedural software components. As described above, these functions may also include any of the steps of Fig. 1. In a preferred embodiment of the invention, these non-inference elements may be used to generate a model of biological activity, as the data is being accumulated and verified.

15 In a preferred embodiment of the invention, system 70 includes a current operating state, which may affect and/or be affected by the inference rules. For example, the number of discovered drug candidates may affect the strictness in which inference rules are applied and/or in which queries are broadened. Preferably, different data matching functions are provided as a function of the system operating state.

20 In a preferred embodiment of the invention, system 70 execute multiple inter-communicating agents simultaneously. Such a configuration may be useful if for example each agent attacks a different aspect of a same problem. In a preferred embodiment of the invention, the agents communicate when one agent discovers knowledge relevant to another agent. Preferably communication uses the KQML standard for communications, using either CLIPS
25 or KIF as an encapsulated knowledge representation language.

In a preferred embodiment of the invention, multi-agent system 70 includes a facilitator agent (or adapter), for example as defined in the FIPA (Foundation for Intelligent Physical Agents) guidelines. This facilitator agent may be used to enable one agent to find a second agent have a specialty of finding a certain type of data. Thus, an agent may function as a data
30 source.

In a preferred embodiment of the invention, system 70 includes a communication coordinator which coordinates the communications of the multiple agents with the external and/or internal databases.

CLAIMS

1. A method of genomic data discovery, comprising:
 - 5 (a) providing a gene data base comprising at least 10 genes;
 - (b) selecting one of said at least 10 genes;
 - (c) discovering knowledge for said selected gene;
 - (d) repeating said (b) and (c) for a plurality of said genes; and
 - (e) repeating said (b)-(d) a plurality of times such that knowledge is discovered
- 10 substantially in parallel for all the selected genes.
2. A method according to claim 1, wherein said (b)-(e) are performed substantially without human intervention.
- 15 3. A method according to any of claims 1-2, comprising evaluating, by a computer and without requiring additional input from an operator, said genes for which knowledge has been discovered.
4. A method according to claim 3, wherein automatically evaluating said genes comprises
- 20 ranking said genes according to their suitability for being drug leads.
5. A method according to any of claims 3-4, comprising deciding on further selecting of said genes in (b), responsive to said evaluation.
- 25 6. A method according to any of claims 1-5, wherein (c) comprises determining data needs for said genes.
7. A method according to claim 6, wherein each of said genes is associated with a scheme and wherein determining data needs comprise analyzing said scheme to determine data needs.
- 30 8. A method according to any of claims 6-7, wherein (c) comprises formulating queries to obtain said needed data.

21. A method according to any of claims 1-15, wherein said plurality of genes comprises at least 80% of said at least 10 genes.
- 5 22. A method of genomic knowledge discovery, comprising:
determining at least one required data element for at least one gene;
querying a plurality of at least 50 databases for said at least one required data element;
receiving responses from said databases; and
analyzing said responses to increase knowledge for said at least one gene.
- 10 23. A method according to claim 22, wherein said at least 50 databases comprise at least 100 databases.
24. A method according to claim 22, wherein said at least 50 databases comprise at least
15 300 databases.
25. A method according to any of claims 22-24, wherein said databases are all queried for a same data value.
- 20 26. A method according to any of claims 22-25, wherein said method is performed , by a computer and without requiring additional input from an operator,.
27. A method of automated knowledge discovery, comprising:
continuously operating a knowledge discovery cycle comprising:
25 querying a database to receive data; and
performing inferences on said data to generate knowledge; and
re-evaluating said inferences when said database is modified
28. A method according to claim 27, wherein said cycle is continuously operated over one
30 week.
29. A method according to claim 27, wherein said cycle is continuously operated over one month.

39. A system according to any of claims 34-36, wherein said at least 10 adapter units comprise at least 300 adapter units for 300 dissimilar data sources.

5 40. A system according to any of claims 34-39, wherein said at least 10 data sources comprise at least 10 data analysis tools.

41. A system according to any of claims 34-39, wherein said at least 10 data sources comprise at least 30 data analysis tools.

10 42. A system according to any of claims 34-41, wherein said adapters are programmed in an interpreted text processing language.

43. A system according to claim 42, wherein said adapters are programmed in language
15 including classes.

44. A system according to claim 42 or claim 43, wherein said adapters are programmed in a Perl-like language.

20 45. A system according to any of claims 34-43, comprising a central adapter registry, wherein, each of said adapter registers its availability in said central registry.

46. A system according to claim 45, wherein said registry is implemented as assertions.

25 47. A system according to claim 45 or claim 46, wherein said adapters register their data provision capabilities in said registry.

48. A system according to any of claims 34-47, wherein said first unit analyses said data requirements to determine selected ones of said data sources to query.

30 49. A method of ranking genes for an application, comprising:
providing a plurality of gene tokens;
applying, by a computer and without requiring additional input from an operator,
application-specific ranking rules to said gene tokens; and

60. A method of genomic information analysis, comprising:
providing a first model of a biological relationship which interrelates a first plurality of
genes or proteins;
providing a second model of a biological relationship which interrelates a second
5 plurality of genes or proteins; and
applying inference rules to said first and second models to infer missing information.
61. A method according to claim 60, wherein said applying comprises determining data
needs for at least one of said models, based on said applied inference rules.
- 10 62. A method according to claim 60 or claim 61, wherein said biological relationships
comprise enzymatic pathways.
63. A method according to any of claims 60-62, wherein said biological relationships are
15 in different species.
64. A method of automated genomic knowledge discovery, comprising:
analyzing a gene token to determine required data;
selecting at least one suitable human expert; and
20 querying the at least one selected human expert for the required data.
65. A method of automated genomic knowledge discovery, comprising:
analyzing a gene token to determine required data; and
generating, by a computer and without requiring additional input from an operator, a
25 work order to a laboratory to generate the required data.

2/4

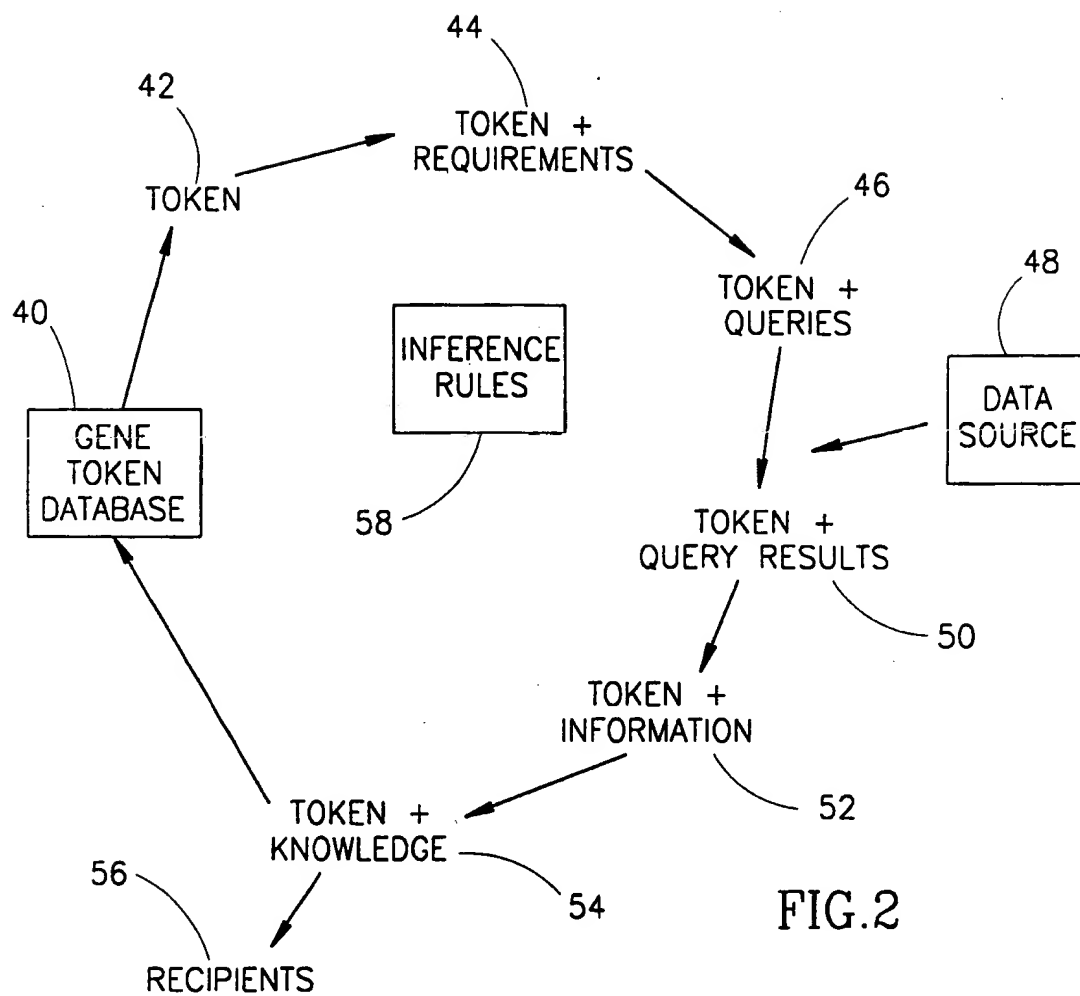


FIG. 2

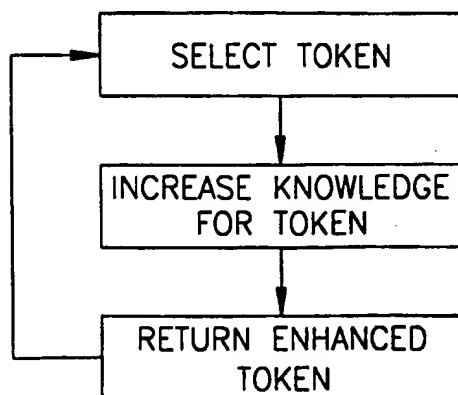


FIG. 3

4/4

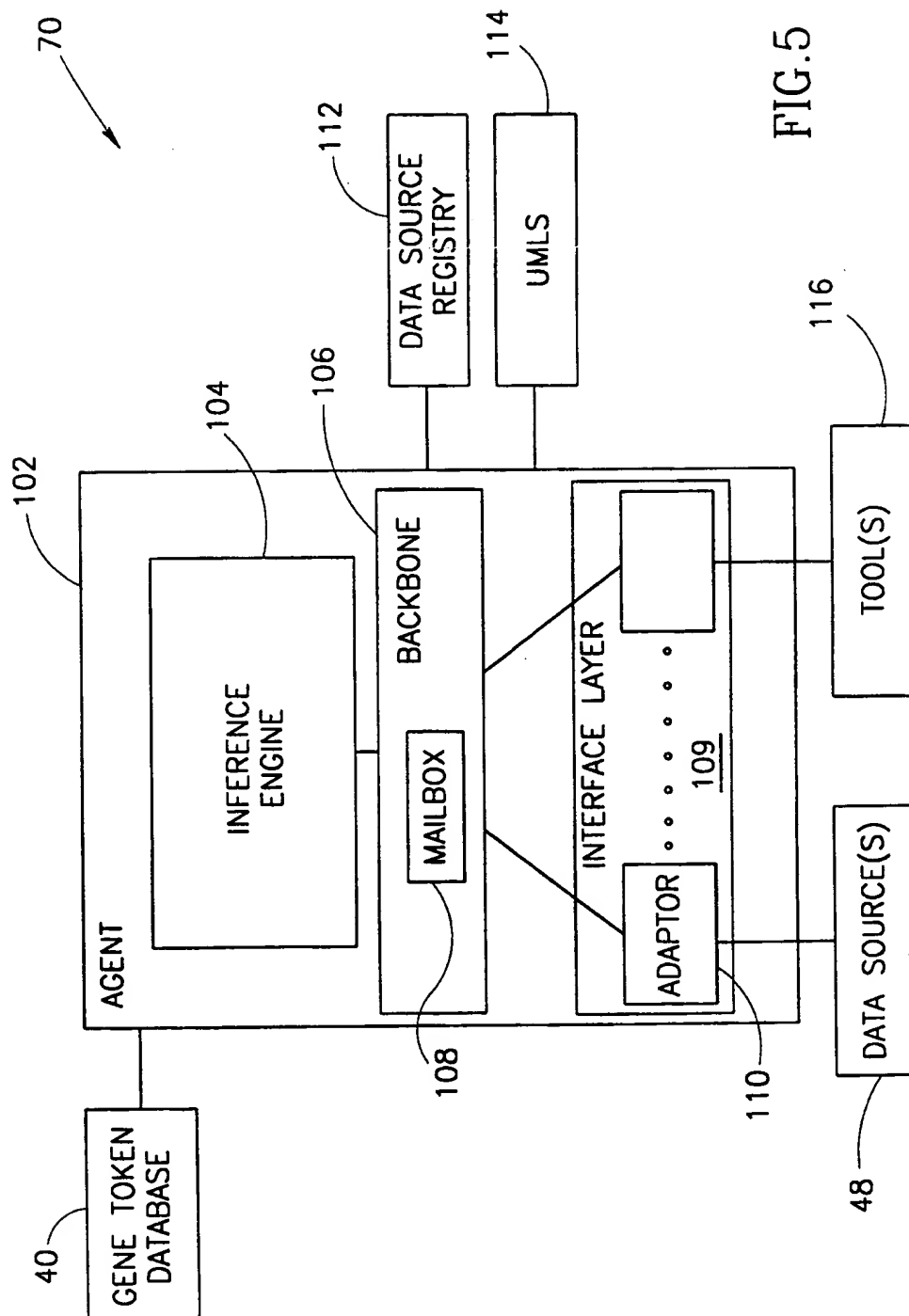


FIG.5